1.0

1.1

1.25

4.5
5.0
5.6

3.6

2.8

3.2

4.0

2.5

2.2

2.0

1.8

1.4

1.6

MICROCOPY RESOLUTION TEST CHART

NATIONAL BUREAU OF STANDARDS-1963-A

ARO 13149.4-M

See 1473
in back

# UNIVERSITY OF MINNESOTA

## SCHOOL OF
## STATISTICS

D D C
FEB 8 1978
D

Simulation Studies on Some Nearest Neighbor Rules

for Statistical Classification.[1]

By

David Aarons and Somesh Das Gupta

University of Minnesota
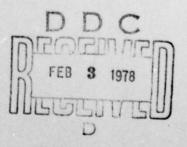Technical Report No. 303

November 1977

DDC

FEB 3 1978

Simulation Studies on Some Nearest Neighbor Rules

for Statistical Classification.[1]


by

David Aarons and Somesh DasGupta

University of Minnesota

1. **Introduction.** The two-population classification problem is to identify a population $\pi_0$ with one of two given populations $\pi_1$ and $\pi_2$ based on observations from these populations on a random vector $X$. We shall consider here $X$ to be univariate. Let $F_i$ be the c.d.f. of $X$ in $\pi_i$ $(i = 0, 1, 2)$. Thus our problem is to test $H_1: F_0 = F_1$ vs. $H_2: F_0 = F_2$. In this paper we have considered some rules which are suggested in the literature when $F_1, F_2$ are not known except that they are continuous. We have studied the performances of the following three rules by simulation.

Let $X_0$, $X_{1i}$ $(i = 1, \ldots, n_1)$, $X_{2i}$ $(i = 1, \ldots, n_2)$ be random observations on $X$ from the populations $\pi_0, \pi_1, \pi_2$, respectively.

**Rule I.** **1-NN (nearest neighbor) Rule:** Measure distances of $X_0$ from $X_{1i}$'s and $X_{2i}$'s and based on these distances classify $X_0$ into the population to which its nearest neighbor belongs.

**Rule II.** **1-RNN (rank nearest neighbor) Rule:** Pool all the observations and order them.

(a) If $X_0$ is the largest or the smallest observation classify $X_0$ into the population of its nearest neighbor (based on ranks).

(b) If both the right-hand and the left-hand nearest neighbor of $X_0$ (denoted by $U_1$ and $V_1$) belong to the same population, classify $X_0$ into that population.

(c) If $U_1$ and $V_1$ belong to different populations classify $X_0$ into $\pi_1$ and $\pi_2$ with probabilities 1/2 and 1/2, respectively. (We call this case a "tie".)

**Rule III.** **2-RNN Rule:** Apply the 1-RNN rule. If a tie occurs, delete the observations corresponding to $U_1$ and $V_1$ and apply the 1-RNN rule again on the remaining observations.

The first rule was suggested and studied by Fix and Hodges (1951, 1953). DasGupta and Lin (1977) proposed the RNN rules and obtained the asymptotic probabilities of misclassification as $n_1$, $n_2 \to \infty$. For a given rule $\delta$, let its PMC under $F_0 = F_1$ be given by

$$\alpha(\delta) = \Pr[\delta \text{ classifies } X_0 \text{ into } \pi_2 \,|\, F_0 = F_1] \,.$$

Let $\alpha_1^*$, $\alpha_2^*$, $\alpha_3^*$ be the asymptotic values of $\alpha$ corresponding to the above rules 1, 2 and 3. Let $f_i$ be the p.d.f. of $F_i$ with respect to Lebesgue measure $(i = 1,2)$ and $p_i = \lim n_i/(n_1 + n_2)$ $(i = 1,2)$ as $\min(n_1, n_2) \to \infty$. It was shown by Fix and Hodges (1951) and DasGupta and Lin (1977) that

$$\alpha_1^* = \alpha_2^* = \int_{-\infty}^{\infty} p_2 f_1(x) f_2(x) dx / \{p_1 f_1(x) + p_2 f_2(x)\}$$

$$\alpha_3^* = \alpha_2^* + \int_{-\infty}^{\infty} \frac{p_1 p_2 f_1(x) f_2(x) \cdot \{p_2 f_2(x) - p_1 f_1(x)\}}{\{p_1 f_1(x) + p_2 f_2(x)\}^3} f_1(x) dx \,.$$

In this paper we have studied the finite-sample performances of these rules by estimating $\alpha$ based on samples from sets of two given populations.

2. **The Experiment.** Different steps of our simulation study are given below.

(i) Two known but different univariate distributions $F_1$ and $F_2$ are chosen.

(ii) Random samples of sizes $n_1$ and $n_2$ from $F_1$ and $F_2$, respectively, are obtained; these samples are called training samples.

(iii) A random sample of size $n_0$ from $F_0 = F_1$ is obtained. We call this a test sample.

(iv) For each observation in the test sample a given classification rule $\delta$ (one of the above three rules) is applied and let $n_{02}$ be the number of the observations in the test sample which are classified by $\delta$

into $F_2$ . Let $\hat{\alpha}(\delta) = n_{02}/n_0$ be the proportion of test samples misclassified into $F_2$ .

(v)  Steps (ii)-(iv) are repeated $r$ times for new training and test samples keeping $n_1$, $n_2$ and $n_0$ fixed.

(vi)  The mean and the standard error of the mean based on $r$ values of $\hat{\alpha}(\delta)$ thus obtained are recorded.

(vii)  Steps (ii)-(vi) are repeated for different values of $n_1$, $n_2$ and $r$ .

(viii) $F_2$ is characterized by a parameter $\theta$ . For different values of $\theta$ steps (i)-(vii) are repeated.

Our choices are given in the following table.

| $F_1$ | $F_2$ | Parameters | $n_1=n_2$ | $n_0$ | $r$ |
|---|---|---|---|---|---|
| $N(0,1)$ | $N(\theta,1)$ | $\theta=0, \pm1, \pm2, 3$ | 25<br>100 | 100<br>400 | 20<br>4 |
| $N(0,1)$ | $N(0,\theta)$ | $\theta=2, 3, 1/2, 1/3$ | 25<br>100 | 100<br>400 | 20<br>4 |
| $e^{-x}$<br>(density) | $\theta e^{-\theta x}$ | $\theta=1, 2, 3, 4,$<br>$1/2, 1/3, 1/4, 1/8$ | 100 | 100 | 20 |
| Cauchy $(0,1)$ | Cauchy $(\theta,1)$ | $\theta=0, \pm1, \pm2, \pm3$ | 25<br>100 | 100<br>400 | 20<br>4 |

Samples are generated by a library subroutine available on the CDC 6400 at the University of Minnesota.

Note 1.  In the following tables "Half" refers to taking one-half the number of ties to count as misclassified and "R-half" refers to resolving the ties by the use of uniform random number generator.

**Note 2.** In some of the following tables EPMC denotes an estimate of the asymptotic PMC $(\alpha_1^* = \alpha_2^*)$ of the 1-NN and 1-RNN rules. These are derived by the method of runs as suggested in Das Gupta and Lin (1977).

3. **Tables**

### Table 3.1

**Proportion of test sample misclassified into $\pi_2$.**

$F_1 = N(0,1)$, $F_2 = N(\theta,1)$; $n_1 = n_2 = 25$, $n_0 = 100$, $r = 20$ .

Optimal (assuming $\theta$ is known and for minimax rule) PMC is

$\Phi(-|\theta|/2)$ .

| $\theta$ \ Rule | 1NN | | RNN | | | 2-RNN | | | Opt. Exp't. |
|---|---|---|---|---|---|---|---|---|---|
| | MEAN | s.e. | | MEAN | s.e. | | MEAN | s.e. | PMC |
| $\theta = 0$ | .479 | .017 | Half | .479 | .013 | Half | .479 | .014 | .500 |
| | | | Rhalf | .485 | .016 | Rhalf | .484 | .015 | |
| $\theta = 1$ | .374 | .018 | Half | .381 | .014 | Half | .343 | .021 | .308 |
| | | | Rhalf | .374 | .016 | Rhalf | .340 | .021 | |
| $\theta = -1$ | .426 | .020 | Half | .426 | .014 | Half | .421 | .025 | .308 |
| | | | Rhalf | .432 | .017 | Rhalf | .425 | .024 | |
| $\theta = 2$ | .195 | .018 | Half | .194 | .018 | Half | .165 | .017 | .159 |
| | | | Rhalf | .196 | .018 | Rhalf | .164 | .018 | |
| $\theta = -2$ | .245 | .020 | Half | .254 | .018 | Half | .258 | .019 | .159 |
| | | | Rhalf | .251 | .018 | Rhalf | .255 | .018 | |
| $\theta = 3$ | .086 | .012 | Half | .089 | .012 | Half | .062 | .010 | .067 |
| | | | Rhalf | .084 | .011 | Rhalf | .061 | .009 | |
| $\theta = -3$ | .105 | .013 | Half | .114 | .012 | Half | .119 | .015 | .067 |
| | | | Rhalf | .113 | .011 | Rhalf | .118 | .015 | |

## Table 3.2

**Proportion of test sample misclassified into $\pi_2$ .**

$F_1 = N(0,1)$, $F_2 = N(\theta,1)$; $n_1 = n_2 = 100$, $n_0 = 400$, $r = 4$ .

| θ \ Rule | 1NN MEAN | s.d. | RNN | MEAN | s.d. | EPMC | 2-RNN | MEAN | s.d. | Opt. Exp't. PMC |
|---|---|---|---|---|---|---|---|---|---|---|
| θ = 0 | .490 | .018 | Half<br>Rhalf | .482<br>.475 | .008<br>.006 | .48 | Half<br>Rhalf | .509<br>.501 | .014<br>.016 | .500 |
| θ = 1 | .415 | .010 | Half<br>Rhalf | .398<br>.404 | .014<br>.024 | .36 | Half<br>Rhalf | .351<br>.358 | .009<br>.024 | .308 |
| θ = -1 | .402 | .010 | Half<br>Rhalf | .394<br>.397 | .007<br>.007 | .38 | Half<br>Rhalf | .347<br>.344 | .025<br>.024 | .308 |
| θ = 2 | .208 | .010 | Half<br>Rhalf | .210<br>.208 | .010<br>.009 | .22 | Half<br>Rhalf | .200<br>.199 | .011<br>.012 | .159 |
| θ = -2 | .209 | .012 | Half<br>Rhalf | .213<br>.215 | .008<br>.009 | .22 | Half<br>Rhalf | .197<br>.200 | .013<br>.014 | .159 |
| θ = 3 | .088 | .011 | Half<br>Rhalf | .083<br>.082 | .009<br>.007 | .10 | Half<br>Rhalf | .065<br>.066 | .005<br>.006 | .007 |
| θ = -3 | .104 | .012 | Half<br>Rhalf | .101<br>.107 | .008<br>.013 | .09 | Half<br>Rhalf | .088<br>.094 | .012<br>.014 | .007 |

### Table 3.3.

**Proportion of test sample misclassified into $\pi_2$ .**

$F_1 = N(0,1)$, $F_2 = N(0,\theta)$; $n_1 = n_2 = 25$, $n_0 = 100$. $r = 20$ .

| $\theta$ \ Rule | 1NN MEAN | s.e. | RNN MEAN | s.e. | 2-RNN MEAN | s.e. |
|---|---|---|---|---|---|---|
| $\theta = 2.0$ | .375 | .009 | Half .394<br>Rhalf .393 | .008<br>.010 | Half .353<br>Rhalf .355 | .014<br>.015 |
| $\theta = 3.0$ | .399 | .014 | Half .346<br>Rhalf .337 | .013<br>.013 | Half .293<br>Rhalf .295 | .019<br>.018 |
| $\theta = .5$ | .417 | .017 | Half .438<br>Rhalf .337 | .015<br>.018 | Half .461<br>Rhalf .460 | .020<br>.021 |
| $\theta - 1/3$ | .359 | .022 | Half .376<br>Rhalf .380 | .018<br>.019 | Half .393<br>Rhalf .391 | .019<br>.019 |

### Table 3.4

**Proportion of test sample misclassified into $\pi_2$ .**

$F_1 = N(0,1)$, $F_2 = N(0,\theta)$; $n_1 = n_2 = 100$, $n_0 = 400$, $r = 4$ .

| $\theta$ \ Rule | 1NN MEAN | s.e. | RNN MEAN | s.e. | EPMC | 2-RNN MEAN | s.e. |
|---|---|---|---|---|---|---|---|
| $\theta = 2.0$ | .435 | .022 | Half .424<br>Rhalf .426 | .022<br>.025 | .36 | Half .395<br>Rhalf .396 | .027<br>.028 |
| $\theta = 3.0$ | .333 | .012 | Half .338<br>Rhalf .336 | .010<br>.011 | .32 | Half .295<br>Rhalf ..296 | .012<br>.011 |
| $\theta = .5$ | .397 | .062 | Half .405<br>Rhalf .407 | .011<br>.013 | .38 | Half .409<br>Rhalf .408 | .006<br>.005 |
| $\theta = 1/3$ | .339 | .021 | Half .352<br>Rhalf .354 | .020<br>.021 | .35 | Half .360<br>Rhalf .361 | .029<br>.030 |

## Table 3.5

Proportion of test sample misclassified into $\pi_2$.

$f_1(x) = e^{-x}$, $f_2(x) = \theta e^{-\theta x}$; $n_1 = n_2 = n_0 = 100$, $r = 4$ .

| $\theta$ \ Rule | 1NN | | RNN | | | EPMC | 2-RNN | |
|---|---|---|---|---|---|---|---|---|
| | MEAN | s.e. | | MEAN | s.e. | | MEAN | s.e. |
| $\theta = 1$ | .508 | .016 | Half .509<br>Rhalf .523 | | .013<br>.016 | .47 | Half .503<br>Rhalf .517 | .013<br>.015 |
| $\theta = 2$ | .442 | .015 | Half .434<br>Rhalf .438 | | .014<br>.017 | .38 | Half .442<br>Rhalf .444 | .016<br>.016 |
| $\theta = 3$ | .402 | .014 | Half .388<br>Rhalf .387 | | .011<br>.011 | .36 | Half .394<br>Rhalf .387 | .013<br>.014 |
| $\theta = 4$ | .335 | .009 | Half .330<br>Rhalf .336 | | .007<br>.008 | .32 | Half .327<br>Rhalf .330 | .009<br>.009 |
| $\theta = .5$ | .453 | .010 | Half .453<br>Rhalf .458 | | .009<br>.013 | .38 | Half .430<br>Rhalf .430 | .010<br>.014 |
| $\theta = 1/3$ | .410 | .011 | Half .395<br>Rhalf .386 | | .008<br>.009 | .36 | Half .346<br>Rhalf .335 | .010<br>.010 |
| $\theta = 1/4$ | .354 | .015 | Half .364<br>Rhalf .372 | | .012<br>.013 | .32 | Half .290<br>Rhalf .292 | .013<br>.013 |
| $\theta = 1/8$ | .247 | .014 | Half .248<br>Rhalf .259 | | .012<br>.014 | .22 | Half .181<br>Rhalf .185 | .011<br>.010 |

## Table 3.6

**Proportion of test sample misclassified into $\pi_2$.**

$F_1$ = Cauchy$(0,1)$, $F_2$ = Cauchy$(\dot\theta,1)$; $n_1 = n_2 = 25$, $n_0 = 100$, $r = 20$ .

| $\theta$ \ Rule | 1NN | | RNN | | | 2-RNN | | |
|---|---|---|---|---|---|---|---|---|
| | MEAN | s.e. | | MEAN | s.e. | | MEAN | s.e. |
| $\theta = 0$ | .473 | .018 | Half | .430 | .015 | Half | .488 | .027 |
| | | | Rhalf | .493 | .018 | Rhalf | .505 | .029 |
| $\theta = 1$ | .406 | .022 | Half | .418 | .022 | Half | .397 | .031 |
| | | | Rhalf | .408 | .025 | Rhalf | .395 | .033 |
| $\theta = -1$ | .398 | .016 | Half | .410 | .012 | Half | .389 | .021 |
| | | | Rhalf | .410 | .013 | Rhalf | .385 | .022 |
| $\theta = 2$ | .288 | .021 | Half | .297 | .021 | Half | 248 | .027 |
| | | | Rhalf | .288 | .021 | Rhalf | .238 | .028 |
| $\theta = -2$ | .247 | .012 | Half | .264 | .012 | Half | .248 | .017 |
| | | | Rhalf | .276 | .015 | Rhalf | .252 | .019 |
| $\theta = 3$ | .161 | .020 | Half | .168 | .017 | Half | .103 | .017 |
| | | | Rhalf | .161 | .018 | Rhalf | .099 | .017 |
| $\theta = -3$ | .153 | .015 | Half | .156 | .013 | Half | .130 | .014 |
| | | | Rhalf | .154 | .013 | Rhalf | .125 | .014 |

## Table 3.7

### Proportion of test sample misclassified into $\pi_2$.

$F_1 = \text{Cauchy}(0,1)$, $F_2 = \text{Cauchy}(\theta,1)$; $n_1 = n_2 = 100$, $n_0 = 400$, $r = 4$ .

| $\theta$ \ Rule | 1NN MEAN | s.e. | RNN | MEAN | s.e. | 2-RNN | MEAN | s.e. |
|---|---|---|---|---|---|---|---|---|
| $\theta = 0$ | .494 | .015 | Half<br>Rhalf | .514<br>.529 | .013<br>.014 | Half<br>Rhalf | .506<br>.512 | .017<br>.021 |
| $\theta = 1$ | .411 | .010 | Half<br>Rhalf | .426<br>.446 | .009<br>.018 | Half<br>Rhalf | .381<br>.390 | .018<br>.017 |
| $\theta = -1$ | .457 | .029 | Half<br>Rhalf | .446<br>.454 | .033<br>.025 | Half<br>Rhalf | .394<br>.393 | .028<br>.025 |
| $\theta = 2$ | .284 | .007 | Half<br>Rhalf | .278<br>.283 | .008<br>.009 | Half<br>Rhalf | .217<br>.219 | .033<br>.024 |
| $\theta = -2$ | .152 | .016 | Half<br>Rhalf | .318<br>.321 | .022<br>.015 | Half<br>Rhalf | .254<br>.257 | .014<br>.010 |
| $\theta = 3$ | .152 | .016 | Half<br>Rhalf | .154<br>.417 | .015<br>.012 | Half<br>Rhalf | .088<br>.087 | .018<br>.014 |
| $\theta = -3$ | .204 | .034 | Half<br>Rhalf | .199<br>.198 | .032<br>.034 | Half<br>Rhalf | .105<br>.103 | .011<br>.012 |

4. <u>Concluding Remarks</u>. For all the three rules considered, it seems that $\hat{\alpha}_1$ has a definite tendency to decrease as $\theta$ moves away (in either direction) from its value under $F_1$.

For small $n_1 = n_2$ there is not any marked difference in performances of these three rules although the 2-RNN rule may be a bit better. However, for large $n_1 = n_2$ the 2-RNN rule seem to have markedly better performance except for the cases $N(0,1)$ vs. $N(0,\theta)$, $\theta < 1$. This report is the first empirical study on the performances of 1NN and RNN rules, although a more detailed study especially on multi-stage RNN rules is called for.

# References

1. Das Gupta, S. and Lin, H. E. (1977). Nearest neighbor rules for statistical classification based on ranks. Tech. Rep. 285, <u>School of Statistics</u>, <u>University of Minnesota</u>, Minneapolis, Minnesota.

2. Fix, E. and Hodges, J. L. (1951). Nonparametric discrimination: Consistency properties. <u>U.S. Air Force School of Aviation Medicine</u>. Report No. 4. Randolph Field, Texas.

3. Fix, E. and Hodges, J. L. (1953). Nonparametric discrimination. Small sample properties, <u>Ibid</u>., Report No. 11.

## REPORT DOCUMENTATION PAGE

**READ INSTRUCTIONS**
**BEFORE COMPLETING FORM**

| | |
|---|---|
| **1. REPORT NUMBER** (14) U MN/DTS/TR-303 2. 30VT ACCESSION NO.<br>Technical Report No. 303 | **3. RECIPIENT'S CATALOG NUMBER** |
| **4. TITLE** *(and Subtitle)*<br>(6) Simulation Studies on Some Nearest Neighbor Rules for Statistical Classification, | **5. TYPE OF REPORT & PERIOD COVERED**<br>(9) Technical Report,<br>**6. PERFORMING ORG. REPORT NUMBER** |
| **7. AUTHOR(s)**<br>(10) Somesh/Das Gupta and David/Aarons | **8. CONTRACT OR GRANT NUMBER(s)**<br>(15) DAAG-29-76-G-0038 |
| **9. PERFORMING ORGANIZATION NAME AND ADDRESS**<br>Department of Theoretical Statistics<br>University of Minnestoa<br>Minneapolis, MN 55455 | **10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS** |
| **11. CONTROLLING OFFICE NAME AND ADDRESS**<br>U.S. Army Research Office<br>P.O. Box 12211<br>Research Triangle Park, NC 27709 | **12. REPORT DATE**<br>(11) November 1977<br>**13. NUMBER OF PAGES**<br>11 (12) 15 P. |
| **14. MONITORING AGENCY NAME & ADDRESS***(if different from Controlling Office)*<br>(18) ARO (19) 13149.4-M | **15. SECURITY CLASS.** *(of this report)*<br>Unclassified |
| | **15a. DECLASSIFICATION/DOWNGRADING SCHEDULE**<br>NA |

**16. DISTRIBUTION STATEMENT** *(of this Report)*

Approved for public release; distribution unlimited.

**17. DISTRIBUTION STATEMENT** *(of the abstract entered in Block 20, if different from Report)*

NA

**18. SUPPLEMENTARY NOTES**

The findings in this report are not to be construed as an official Department of Army position, unless so designated by other authorized documents.

**19. KEY WORDS** *(Continue on reverse side if necessary and identify by block number)*

Classification, nearest neighbor rule, rank nearest nearest neighbor rule, simulation, probability of misclassification.

**20. ABSTRACT** *(Continue on reverse side if necessary and identify by block number)*

Abstract: The performances of the distance nearest neighbor rule, one-stage rank nearest neighbor rule and two-stage rank nearest neighbor rule for the two-population statistical classification problem are studied through simulation processes. Normal distributions with different mean, normal distributions with different variances, exponential distributions and Cauchy distributions are considered.

**DD** FORM 1473 EDITION OF 1 NOV 65 IS OBSOLETE

410 022 JeB